

Worker





in "big data" systems















Basic architecture

Clients submit applications to the cluster manager Cluster manager assigns cluster resources to applications Each Worker launches containers for each application Driver containers run main method of user program Executor containers run actual computation

Examples of cluster manager: YARN, Mesos *Examples* of computing frameworks: Hadoop MapReduce, Spark

12

Two levels of scheduling

Cluster-level: Cluster manager assigns resources to applications Application-level: Driver assigns *tasks* to run on executors A task is a unit of execution that operates on one *partition*

Some advantages:

Applications need not be concerned with resource fairness Cluster manager need not be concerned with individual tasks Easy to implement priorities and preemption

13



Case Study: MapReduce

(Data-parallel programming at scale)

Application: Word count

Locally: tokenize and put words in a hash map

How do you parallelize this?

Split document by half

Build two hash maps, one for each half

Merge the two hash maps (by key)

16

























MapReduce

Partition dataset into many chunks

Map stage: Each node processes one or more chunks locally

Reduce stage: Each node fetches and merges partial results from all other nodes

29

31

MapReduce: Word count



// key = document name
// value = document contents
for each word w in value:
 emit (w, 1)

reduce(key, values):

// key = the word
// values = number of occurrences of that word
count = sum(values)
emit (key, count)









Brainstorm: Top K

Find the largest K values from a set of numbers

How would you express this as a distributed application? In particular, what would map and reduce phases look like?

Hint: use a heap...

Brainstorm: Top K

Assuming that a set of K integers fit in memory...

Key idea...

35

Map phase: everyone maintains a heap of K elements

Reduce phase: merge the heaps until you're left with one

36

38

40

Brainstorm: Top K

Problem: What are the keys and values here?

No notion of key here, just assign the same key to all the values (e.g. key = 1)

Map task 1: $[10, 5, 3, 700, 18, 4] \rightarrow (1, heap(700, 18, 10))$ Map task 2: $[16, 4, 523, 100, 88] \rightarrow (1, heap(523, 100, 88))$ Map task 3: $[3, 3, 3, 3, 300, 3] \rightarrow (1, heap(300, 3, 3))$ Map task 4: $[8, 15, 20015, 89] \rightarrow (1, heap(20015, 89, 15))$

Then all the heaps will go to a single reducer responsible for the key 1 This works, but clearly not scalable...

37

Brainstorm: Top K

Idea: Use X different keys to balance load (e.g. X = 2 here) Map task 1: [10, 5, 3, 700, 18, 4] → (1, heap(700, 18, 10)) Map task 2: [16, 4, 523, 100, 88] → (1, heap(523, 100, 88)) Map task 3: [3, 3, 3, 3, 300, 3] → (2, heap(300, 3, 3)) Map task 4: [8, 15, 20015, 89] → (2, heap(20015, 89, 15))

Then all the heaps will (hopefully) go to X different reducers Rinse and repeat (*what's the runtime complexity?*)

Monday 3/25

Stream processing

Application: Word Count

SELECT count(word) FROM data GROUP BY word

cat data.txt

| tr -s '[[:punct:][:space:]]' '\n'

| sort | uniq -c

42

44

Using partial aggregation

- 1. Compute word counts from individual files
- 2. Then merge intermediate output
- 3. Compute word count on merged outputs

Using partial aggregation

- 1. In parallel, send to worker:
 - Compute word counts from individual files
 - Collect result, wait until all finished
- 2. Then merge intermediate output
- 3. Compute word count on merged intermediates

MapReduce: Programming Interface

map(key, value) -> list(<k', v'>)

 Apply function to (key, value) pair and produces set of intermediate pairs

43

reduce(key, list<value>) -> <k', v'>

- Applies aggregation function to values
- Outputs result

MapReduce: Programming Interface

map(key, value):
 for each word w in value:
 EmitIntermediate(w, "1");

```
reduce(key, list(values):
    int result = 0;
    for each v in values:
        result += ParseInt(v);
    Emit(AsString(result));
```

46

48

MapReduce: Optimizations

combine(list<key, value>) -> list<k,v>

- Perform partial aggregation on mapper node: $\langle \text{the, 1} \rangle, \langle \text{the, 1} \rangle \rightarrow \langle \text{the, 3} \rangle$
- reduce() should be commutative and associative

partition(key, int) -> int

- Need to aggregate intermediate vals with same key

45

47

- Given n partitions, map key to partition $0 \le i < n$
- Typically via hash(key) mod n

Fault Tolerance in MapReduce



Map worker writes intermediate output to local disk, separated by partitioning. Once completed, tells master node.

 Reduce worker told of location of map task outputs, pulls their partition's data from each mapper, execute function across data

Note:

- "All-to-all" shuffle b/w mappers and reducers
- Written to disk ("materialized") b/w each stage

Fault Tolerance in MapReduce

- Master node monitors state of system
 - If master failures, job aborts and client notified
- Map worker failure
 - Both in-progress/completed tasks marked as idle
 - Reduce workers notified when map task is re-executed on another map worker
- Reducer worker failure
 - In-progress tasks are reset to idle (and re-executed)
 - Completed tasks had been written to global file system

Straggler Mitigation in MapReduce



- Tail latency means some workers finish late
- For slow map tasks, execute in parallel on second map worker as "backup", race to complete task